

CLAIMS

WE CLAIM:

1. A method for classifying electronically posted documents, the method comprising:
2. receiving a first document and a second document;
3. generating a first metadata summary corresponding to said first document and a
4. second metadata summary corresponding to the second document, wherein the first metadata
5. summary includes a first summary sub-tree and the second metadata summary includes a second
6. summary sub-tree;

7. comparing the structure of the first summary sub-tree with the structure of the second
8. summary sub-tree; and

9. identifying the first and second documents as distinct if the structures of the first and
10. second summary sub-trees are not equivalent.

11. 2. The method of claim 1, wherein the first summary sub-tree includes at least one
12. attribute having a first attribute value, and wherein the second summary sub-tree includes at least one
13. attribute having a second attribute value, the method further comprising:

14. comparing, for each of the at least one attributes, the first and second attribute values;
15. and

16. identifying the first and second documents as distinct if the attribute values of the first
17. and second summary sub-trees are not equivalent.

18. 3. The method of claim 1, wherein the first summary sub-tree includes text content, and
19. wherein the second summary sub-tree includes text content, the method further comprising:

20. comparing the text content included within the first and second summary sub-trees;
21. and

22. identifying the first and second documents as distinct if the text content of the first
23. and second summary sub-trees are not equivalent.

4. The method of claim 2, wherein the first summary sub-tree further includes text content, and wherein the second summary sub-tree includes text content, the method further comprising:

comparing the text content included within the first and second summary sub-trees;

identifying the first and second documents as distinct if the text content included within the first and second summary sub-trees are not equivalent.

5. The method of claim 4, further comprising identifying the first and second documents as duplicates if the text content within the first and second summary sub-trees are equivalent.

6. The method of claim 5, further comprising removing the second metadata summary from the first summary group if the structures of the first and second summary sub-trees are equivalent and if the first summary value is equivalent to the second summary value for each of the at least one attributes.

7. The method of claim 1, further comprising:

defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;

a second row corresponding to the second metadata summary;

a first column corresponding to the first metadata summary; and

6 a second column corresponding to the second metadata summary, wherein the
7 process of identifying the first and second documents as distinct if the structures of the first and
8 second summary sub-trees are not equivalent comprises storing a zero binary value in the first row
9 and second column position of the equivalence metadata summary.

8. The method of claim 2, further comprising:

defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;

a second row corresponding to the second metadata summary;

a first column corresponding to the first metadata summary; and
a second column corresponding to the second metadata summary, wherein the
process of identifying the first and second documents as distinct if the attribute values of the first and
second summary sub-trees are not equivalent comprises storing a zero binary value in the first row
and second column position of the equivalence metadata summary.

9. The method of claim 3, further comprising:

defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;

a second row corresponding to the second metadata summary;

a first column corresponding to the first metadata summary; and

a second column corresponding to the second metadata summary, wherein the process of identifying the first and second documents as distinct if the text content of the first and second summary sub-trees are not equivalent comprises storing a zero binary value in the first row and second column position of the equivalence metadata summary.

10. A method for classifying electronically posted documents, the method comprising:

receiving a plurality of documents;

generating a respective plurality of metadata summaries corresponding to the plurality

4 of received documents;

grouping a first subset of the respective plurality of metadata summaries into a first

6 summary group, the first summary group comprising summaries having a first mime-type

7 designation;

selecting a first metadata summary and a second metadata summary from the first

9 summary group, wherein the first metadata summary includes a first summary sub-tree and the

10 second metadata summary includes a second summary sub-tree;

11 comparing the structure of the first summary sub-tree with the structure of the second

12 summary sub-tree; and

13 identifying the first and second documents as distinct if the structures of the first and

14 second summary sub-trees are not equivalent.

1 11. The method of claim 10, wherein grouping further comprises grouping a second
2 subset of the respective metadata summaries into a second summary group, the second summary
3 group comprising summaries having a second mime-type designation.

1 12. A system for classifying electronically posted documents, the system comprising:
2 a metadata parser module coupled to receive electronically posted documents, the
3 metadata parser configured to output respective metadata summaries, wherein each respective
4 metadata summary comprises one or more sub-trees structures, one or more attributes, and content
5 text;

6 a summary repository coupled to receive and store the respective metadata
7 summaries; and

8 a summary consolidator coupled to the summary repository, the summary
9 consolidator configured to delete duplicate metadata summaries from the summary repository.

10 13. The system of claim 12, wherein the summary consolidator comprises:
11 a sub-tree comparator configured to compare one or more sub-tree structures of the
12 retrieved metadata summaries;

13 an attribute comparator configured to compare the attribute values of the retrieved
14 metadata summaries; and

15 a text comparator configured to compare the text content included within the retrieved
16 metadata summaries.

17 14. The system of claim 13, wherein the sub-tree comparator is configured to compare the
18 metadata portion of the metadata summary.

19 15. The system of claim 13, wherein the attribute comparator is configured to compare
20 the attribute values included within the metadata portion of the metadata summary.

21 16. The system of claim 13, wherein the text comparator is configured to compare the
22 text content included within the metadata portion of the metadata summary.

1 17. A program product for use in a computer system that executes program steps recorded
2 in a computer-readable media to perform a method for classifying electronically posted documents,
3 the program product comprising:
4 a recordable media;
5 a program of computer-readable instructions executable by the computer system to
6 perform processes comprising:
7 receiving a first document and a second document;
8 generating a first metadata summary corresponding to said first document and
9 a second metadata summary corresponding to the second document, wherein the first metadata
10 summary includes a first summary sub-tree and the second metadata summary includes a second
11 summary sub-tree;
12 comparing the structure of the first summary sub-tree with the structure of the
13 second summary sub-tree; and
14 identifying the first and second documents as distinct if the structures of the
15 first and second summary sub-trees are not equivalent.

1 18. The program product of claim 17, wherein the first summary sub-tree includes at least
2 one attribute having a first attribute value, and wherein the second summary sub-tree includes at least
3 one attribute having a second attribute value, the program product method further comprising the
4 processes of:
5 comparing, for each of the at least one attributes, the first and second attribute values;
6 and
7 identifying the first and second documents as distinct if the attribute values of the first
8 and second summary sub-trees are not equivalent.

1 19. The program product of claim 18, wherein the first summary sub-tree includes text
2 content, and wherein the second summary sub-tree includes text content, the program product further
3 comprising the processes of:
4 comparing the text content included within the first and second summary sub-trees;
5 and

identifying the first and second documents as distinct if the text content of the first and second summary sub-trees are not equivalent!.

20. The program product of claim 19, further comprising the method step of identifying the first and second documents as duplicates if the text content within the first and second summary sub-trees are equivalent.

21. The program product of claim 20, further comprising the process of removing the second metadata summary from the first summary group.

22. The program product of claim 21, further comprising the processes of:
defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;

a second row corresponding to the second metadata summary;

a first column corresponding to the first metadata summary; and

a second column corresponding to the second metadata summary, wherein the

process of identifying the first and second documents as distinct if the text content of the first and second summary sub-trees are not equivalent comprises storing a zero binary value in the first row and second column position of the equivalence metadata summary.

23. The method of claim 18, further comprising the processes of:
defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;

a second row corresponding to the second metadata sum

a first column corresponding to the first metadata summary

a second column corresponding to the second metadata summary.

the first and second documents as distinct if the attribute values of the first and trees are not equivalent comprises storing a zero binary value in the first row position of the equivalence metadata summary.

24. The method of claim 19, further comprising the processes of:
defining a first equivalence metadata table comprising:

a first row corresponding to the first metadata summary;
a second row corresponding to the second metadata summary;
a first column corresponding to the first metadata summary; and
a second column corresponding to the second metadata summary, wherein the process of identifying the first and second documents as distinct if the text content of the first and second summary sub-trees are not equivalent comprises storing a zero binary value in the first row and second column position of the equivalence metadata summary.

25. A program product for use in a computer system that executes program steps recorded in a computer-readable media to perform a method for classifying electronically posted documents, the program product comprising:

a recordable media;

a program of computer-readable instructions executable by the computer system to perform method steps comprising:

- receiving a plurality of documents;
- generating a respective plurality of metadata summaries corresponding to the plurality of received documents;
- grouping a first subset of the respective plurality of metadata summaries into a first summary group, the first summary group comprising summaries having a first mime-type designation;
- selecting a first metadata summary and a second metadata summary from the first summary group, wherein the first metadata summary includes a first summary sub-tree and the second metadata summary includes a second summary sub-tree;
- comparing the structure of the first summary sub-tree with the structure of the second summary sub-tree; and
- identifying the first and second documents as distinct if the structures of the first and second summary sub-trees are not equivalent.

26. The program product of claim 25, wherein the step of grouping further comprises the step of grouping a second subset of the respective metadata summaries into a second summary group, the second summary group comprising summaries having a second mime-type designation.